# BY10 Alternatives Analysis
## for the
## Lattice QCD Computing Project Extension
## (LQCD-ext)

*Operated at*
Brookhaven National Laboratory
Fermi National Accelerator Laboratory
Thomas Jefferson National Accelerator Facility

*for the*
U.S. Department of Energy
Office of Science
Offices of High Energy and Nuclear Physics

Version 1.4

Revision Date
August 6, 2009

PREPARED BY:
Chip Watson, JLab

CONCURRENCE:

_____          ____August 6, 2009_____
William N. Boroski                        Date
LQCD Contract Project Manager

**Lattice QCD Computing Project Extension (LQCD-ext)**
**Change Log: Alternatives Analysis for FY10 Submission**

| Revision No. | Description | Effective Date |
|---|---|---|
| 1.0 | Document created for LQCD-ext project. | Dec 18, 2008 |
| 1.1 | Updated for CD-1. | April 13, 2009 |
| 1.2 | Miscellaneous minor corrections. | April 15, 2009 |
| 1.3 | Revised section 3 to state that alternative 1 is the CD-1 preferred down-select option. | June 29, 2009 |
| 1.4 | Adjusted for revised CD2/CD3 budget | August 6, 2009 |
| | | |

# Table of Contents

# 1    Introduction

This document presents the BY10-BY11 analysis of alternatives for the obtaining the computational capacity needed for the US Lattice QCD effort within High Energy Physics (HEP) and Nuclear Physics (NP) by the SC Lattice QCD Computing Extension Project (LQCD-ext). This analysis is updated at least annually to capture decisions taken during the life of the project, and to examine options for the next year of the project. The technical managers of the project are also continuously tracking market developments through interactions with computer and chip vendors, through trade journals and online resources, and through computing conferences. This tracking allows unexpected changes to be incorporated into the project execution in a timely fashion.

Alternatives herein are constrained to approximately fit within the current budget guidance of the project, ~$3.5M / year for the five years of the project (FY10-FY14). This constraint provides adequate funding to meet the basic requirements of the field for enhanced computational capacity, under the assumption of expanding resources at ANL and ORNL already planned by the Office of Science (SC), and under the assumption that a reasonable fraction of those resources are ultimately allocated to Lattice QCD.

All alternatives assume the continued operation of the existing resources from the FY06-FY09 LQCD Facilities Project until those resources reach end of life, i.e., until each resource is no longer cost effective to operate, about 4 years for clusters. In FY10 this will be an aggregate resource of about 18 TFlop/s sustained on LQCD benchmarks. The allocated project cost of operating these existing clusters in FY2010 is approximately $0.9M (for the three sites combined). Replacing (and running) the computational capacity represented by existing resources cannot be done for less than its operating cost.


# 2    FY10-FY11 Goals

The proposed objective for the late FY10 plus early FY11 procurements is to assemble new computational resources that sustain a total of 23 teraflop/s for production lattice QCD calculations. Sustained performance is defined as the average of single precision DWF and single precision improved staggered actions at a scale appropriate for analysis jobs, or 0.1 to 0.5 teraflop/s. "Linpack" or "peak" performance metrics are not considered, as lattice QCD codes uniquely stress computer systems, and their performance does not uniformly track either Linpack or peak performance metrics across different architectures.

The goal for FY10 is to install this new resource by Sept 30, 2010 and release it to full production by Jan 2, 2010. The second phase for this machine, in FY11, will be ordered Q1 FY11, delivered Q2 FY11, and in full production by March 1, 2011. These two procurements are timed to take advantage of a next generation server chip from Intel (major step in price/performance).

A second goal for FY10 is to deliver from all resources 18 teraflop/s-years of running in FY10 on the combination of existing resources, equivalent to an uptime of >90%. The LQCD-ext project defines 1 teraflop/s-year of integrated running time as the result of running a system capable of 1 teraflop/s sustained computing performance on LQCD applications (specifically, the average of the single precision DWF and single precision improved staggered action inverters) for 8000 hours.

Beyond FY11, the objective is to take advantage of the improvements in technology implied by Moore's law, as well as the specific nature of LQCD calculations, to deploy a series of increasingly powerful resources for science.

## 3   Alternatives

The following sections summarize the alternative technologies considered to achieve the stated performance goals of this investment for FY10 and FY11.  Alternative 1, a single optimal cluster deployed in Q4 2010 and Q1 2011 is the Critical Decision 1 (CD-1) preferred down-select option.  This option satisfies FY10-11 goals at the lowest incremental and total lifecycle cost.

### 3.1.   Alternative 1: A Single Optimal Cluster Deployed in Q4 2010 and Q1 2011.

*Deploy an application-optimized cluster in the fourth calendar quarter of 2010 and the first quarter of 2011 to sustain at least 28 teraflop/s on the LQCD single precision analysis benchmarks.*

At the end of FY08, a computer cluster comparable to what is required cost $0.23/MFlop/s to procure.  Based upon current market trends (conservatively 3% average improvement per month in price/performance, equivalent to a 24 month doubling in performance per dollar) this alternative should cost about $0.13 per sustained megaflop/s in late FY2010, thus $1.5M for compute hardware to which we must add $0.24M labor, site prep and overhead.  This cluster would be doubled early in the next fiscal year, yielding a total cost of $3.2M.  Estimated annual operating costs are approximately 10% of the original hardware purchase cost, or $0.32M/yr over an expected 4 year lifetime.

The incremental lifecycle cost of this alternative is estimated as follows:
- Procure and install 23 TF in FY10-FY11 ($3.2M).
- Operations at 10%/yr ($0.32M/yr, $1.3M)
- Incremental lifecycle cost: $4.5M

The total lifecycle cost for the project, with this alternative, is estimated as follows:
- Operate existing computing resources at BNL, FNAL and JLab in FY10: $0.9M
- Operate existing resources at FNAL and JLab in FY11: $0.8M
- Operate existing resources at FNAL in FY12: $0.5M
- FY10-12 Operations plus incremental lifecycle cost for Alternative 1: $6.7M

Risk adjustment: This procurement will be made using a fixed price contract based upon allocated funds, so there is minimal cost risk.  The slower than Moore's Law extrapolation of cost per teraflop/s is both conservative and only assumed for up to 6 months ahead of the average planned procurement dates so as to minimize the risk to the project's goals.  This is because Moore's Law is not continuous but has discrete steps with 3-6 month "flat spots".  Another way to state this is that there is a planned 25% contingency on the minimum performance goal compared to an optimistic scenario.  If the cluster does outperform the minimum stated goal, the science campaigns can easily exploit the additional resource.  In summary there is minimal cost risk, and modest uncertainty in how much the investment may exceed goals.  It should be noted that this consideration of risk uniformly applies to all alternatives, and so has no impact upon the alternatives analysis outcome.

Justification for the expectation of increases in cluster performance/dollar: (1) 6 core processors will be commoditized; (2) DDR-3 1600 memory will be commoditized; (3) processors will increase the number of memory channels supported; (4) QDR Infiniband will be further commoditized.

## 3.2. Alternative 2: Traditional Supercomputers

*Expand the major DOE Supercomputer Centers, National Energy Research Scientific Computing Center (NERSC, Lawrence Berkeley Lab) and the Center for Computational Sciences (CCS, Oak Ridge National Lab) to meet the needs of the QCD physics calculations (an additional 23 teraflop/s roughly in the timescale of quarter 4 of 2010), while continuing to operate the existing systems (FY07-FY09 project clusters and the QCDOC).*

To estimate the price/performance of general use commercial supercomputers, information from a recent supercomputer is used. ORNL's recent $77M upgrade to Jaguar added 1.38 petaflop/s (Linpack) and 150 thousand cores to that machine. Since LQCD achieves about 0.9 GFlop/s per core sustained, this upgrade is about 135 TFlop/s sustained on LQCD, or $0.57/Mflop/s. Assuming Moore's Law, a machine (upgrade) in the desired timeframe (perhaps an XT-6) would cost $0.22/Mflop/s. Note that this is somewhat speculative in that nothing is known about Cray's or other supercomputer vendor's planned release cycles for more cost effective machines, and a machine of this price performance might not be available in this timeframe. In favor of this assumption however is the existence today of faster AMD chips, and a reasonable estimated timetable for yet faster AMD (more cores) or more powerful Intel chips (more cores, more performance per core) that could power a faster Cray.

The incremental lifecycle cost of this alternative is estimated as follows:
- Procure 23 TF in Q4 2010 ($5.1M).
- Operations at 10%/year: ($0.51M/year, $2.0M).
- Incremental lifecycle cost: $5.1M + $2.0M = $7.1M

The total lifecycle cost for the investment, with this alternative, is estimated as follows:
- Operate existing computing resources at BNL, FNAL and JLab in FY10: $0.9M
- Operate existing resources at FNAL and JLab in FY11: $0.8M
- Operate existing resources at FNAL in FY12: $0.5M
- FY10-12 operations plus incremental lifecycle cost for Alternative 2: $9.3M

Analysis: This option falls considerably outside of the investment budget envelope in FY10. High-end clusters, such as the XT-5 at ORNL, which are configured to support a wide range of high end computing applications, contain components (cost factors) which are not needed or are not cost effective for lattice QCD analysis jobs, including higher performance message passing fabrics, higher performance and capacity local disks, and much larger memories. In addition, they are usually integrated with much higher performance file services than required by the LQCD community. (Note: these larger machines are, however, cost competitive for the most demanding LQCD lattice configuration generation jobs, which will be run on supercomputers, but are not covered in this project).

## 3.3 Alternative 3: BlueGene/P Supercomputer

*Purchase (or expand) an IBM BlueGene/\* supercomputer of sufficient size to sustain (an additional) 23 teraflop/s in Q4 2010.*

Purchase a BlueGene/\* commercial supercomputer (BlueGene/P, since the BG/Q is not expected to be available), locate it at one of the labs doing the experiments, and dedicate it to LQCD physics calculations. Or, add additional racks at an existing DOE BG/P site (e.g. ANL). For BlueGene/P, estimates of cost are about $1.25M / rack (assuming 50% discount for age, see below), which contains 1024 processors, or 4096 cores. The BG/P sustains no more than 25% of peak on lattice QCD, or about 0.7 GFlop/s/core. Thus, estimated cost will be $0.43 per Mflop/s. A custom machine or rack might be configured somewhat more cheaply ($0.40). A 28 Tflop/s machine, or partition of a larger machine, will therefore cost $12M.

Note: The BlueGene/Q is expected to be much more cost effective than the BG/P, perhaps by a factor of 8. Release date for this machine is of course not yet known, but mid 2011 is an estimate. Thus, this machine will arrive too late for the first procurement of this project, but might be competitive for the second, FY2012. Evolution of this future machine will be tracked each year just ahead of procurements to evaluate its suitability for inclusion in a benchmarking process.

The support contract for a BG/\* would be included in the procurement price for the first year, and (based upon historical data) would be approximately 8% of the purchase price in subsequent years. Other annual operating costs (staff) are estimated at approximately 2% of the original purchase price.

The incremental 4 year lifecycle cost of this alternative is estimated as follows:
- Procure 23 TF in FY10 ($9M)
- Support contract free in FY10, 8% thereafter: $0.7M/yr
- Operations at 2%/yr: = $0.18M
- Incremental lifecycle cost = $9M + $.2M + $0.9M/yr*3 yrs) = $12M.

The total lifecycle cost for the investment, with this alternative, is estimated as follows:
- Operate existing computing resources at BNL, FNAL and JLab in FY10: $0.9M
- Operate existing resources at FNAL and JLab in FY11: $0.8M
- Operate existing resources at FNAL in FY12: $0.5M
- FY10-12 operations plus incremental lifecycle cost for Alternative 3: $14M

While normally more cost effective than traditional supercomputers when first released, the BlueGene/\* line only takes advantage of Moore's Law about every 3 years, and would still be a year away from the next major gain. The price used here is 50% below initial release prices, showing the effect of selling an "older" machine. Even if sold at 25% of the initial price, this alternative is also too expensive, and if constrained to the project budget, would yield less than half of the incremental capacity that clusters will deliver, at almost twice the average annual cost over its lifetime. Stated another way, it would produce less science per dollar spent.

### 3.4. Alternative 4: Status Quo (no additional deployment in FY10 or FY11)

*Continue to operate the existing project clusters deployed at FNAL and JLab.*

Clusters deployed in FY2007-FY2009 would provide 80% of planned integrated running in FY2010, decreasing to 30% in August 2011 compared with the preferred alternative of a cluster per year. This alternative is included only for completeness and would not be capable of providing the necessary computational capacity to achieve the scientific goals of this project. The cost of this choice is $0.9M in FY2010 to operate the existing facilities. The incremental cost of this alternative (new investment) is $0.

The total project cost with this alternative is estimated as follows:
- Operate existing computing resources at BNL, FNAL and JLab in FY10: $0.9M
- Operate existing resources at FNAL and JLab in FY11: $0.8M
- Operate existing resources at FNAL in FY12: $0.5M
- FY10-12 operations plus incremental lifecycle cost for Alternative 4: $2.2M

### 3.5. Other Alternatives

Other alternatives may be relevant for future iterations of this document. These were not considered for detailed analysis at this time, as their current state of maturity was not deemed sufficient. Each of these alternatives functions as a co-processor with limited memory size, and so could most likely only be used to accelerate floating point intensive kernels. Speed-up will be limited by Amdahl's Law (serial code is not accelerated). The alternatives include:

- IBM/Sony Cell processor based systems: memory bandwidth considerations do not make the current processors sufficiently promising compared to latest generation general purposes processors, but as follow on commercial products emerge additional analysis should be performed.

- General Purpose Graphics Processing Units: Potentially more interesting is the emergence of high performance and programmable GPU's (graphics processing units) with highly parallel vector processing capabilities. These may have sufficient memory bandwidth to allow their use in a mixed CPU/GPU cluster system.

- Other novel accelerators exist, such as ClearSpeed's accelerator board. Like the Cell processor, this chip does not appear to have enough memory bandwidth to sustain the high performance necessary for lattice QCD.

## 4  Net Present Value Considerations

Alternatives 1, 2 and 3 above have the same net present value considerations in that all alternatives have the same computational capacity increment for FY2010, and so yield the same stream of benefits. Thus NPV calculations will have no impact upon the selection of the best alternative, and will give the same rank ordering of the alternatives as does comparing costs.  For completeness, however, we include NPV calculations for alternatives 1, 2, and 3 based on both cost avoidance and on estimated benefits (see below).   Note that the "Baseline: status quo" alternative is rejected in that it does not meet stated goals for capacity.

For both the cost avoidance and estimated benefits NPV calculations, we assume that the hardware becomes operational for production use at the end of Jan 2011 (average of two phases). First, we compare the cost of the chosen alternative (Alternative 1: optimal clusters) to the next most cost effective alternative (Alternative 2: Cray XT-(6)) and use the difference in each year (acquisition and operations costs in the first two years, and operations costs in the subsequent years) as our cost avoidance.  In the second analysis, we use the estimated benefits of $14M for the FY10 investment from the discussion of Return on Investment (see Section 5).

The following table shows the NPV calculated using cost avoidance.

| Cost Avoidance NPV | | | | | | |
|---|---|---|---|---|---|---|
| Discount Rate | 4.3% | | | | | |
| Cost Avoidance | FY10 | FY11 | FY12 | FY13 | FY14 | NPV |
| 1: Cluster | $1.96 | $0.20 | $0.20 | $0.20 | $0.20 | $2.74 |
| 2: Traditional SC | $0.00 | $0.00 | $0.00 | $0.00 | $0.00 | $0.00 |
| 3: BlueGene/P | -$4.14 | $0.32 | -$0.41 | -$0.41 | -$0.41 | -$5.06 |

Table 2 below shows the NPV calculated using an assumed benefit of $12M over a four year lifetime of the system purchased in FY10.  See "Return on Investment" below for discussion of this assumed benefit.  For the NPV calculation, the $12M benefit is spread evenly across the 4 years, with the new systems becoming operational at the end of July 2010.

| Estimated Benefit ($M) | | | | | | |
|---|---|---|---|---|---|---|
| | FY10 Net | FY11 Net | FY12 Net | FY13 Net | FY14 Net | NPV |
| 1: Cluster | -$4.01 | $2.89 | $3.19 | $3.69 | $3.69 | $7.751 |
| 2: Traditional SC | -$5.96 | $2.70 | $3.00 | $3.50 | $3.50 | $5.201 |
| 3: BlueGene/P | -$10.10 | $3.02 | $2.59 | $3.09 | $3.09 | $0.478 |

## 5  Return on Investment

ROI is difficult to quantify, but the following discussion gives an order of magnitude estimate of the benefit of this investment to the HEP and NP programs.

This investment provides two classes of benefits to the High-Energy Physics (HEP) and Nuclear Physics (NP) programs of the DOE's Office of Science (SC). The first class is the direct enhancements to the science itself: these theoretical calculations are important on their own and will lead to new discoveries.  The second is that these calculations are in some cases essential to the cost effective exploitation of much more expensive experiments built and operated by the two program offices.  In the FY07 Operating Plan, the total HEP and NP programs in SC were funded at $752M and $432M, respectively.  Further, both fields of science receive substantial, though smaller, grants from the National Science Foundation. This should be compared to the budget of this project, ~$3.5M/year.  In HEP, roughly 30% of the Tevatron program at Fermilab has a direct interplay with lattice QCD calculations. The whole suite of measurements and calculations are worth much more together than in isolation, so one must conclude that the return on investment for HEP is at least five-fold, and likely higher. In NP, the situation is much the same. A significant development at BNL's Relativistic Heavy-Ion Collider (RHIC) is to search for the critical point of the QCD phase transition. Lattice QCD calculations indicate that this search is within RHIC's reach; RHIC would not proceed without this guidance. At Jefferson Lab a key motivation for the upgraded accelerator is the search for hybrid mesons and gluonic excitations, states whose theoretical foundation rests on lattice QCD. One concludes again that the return on investment for NP is at least five-fold, likely more. With such high rates of return, it is safe to view the calculations as necessary for the DOE to do a sensible deployment of the experiments. But one should then ask whether other computing facilities could do the job. Indeed, all of the experiments in question have computing budgets that rival or surpass this project. However, their communications networks are ill-suited to the data structures of lattice QCD, with a mismatch in efficiency of nearly a factor of 10. In the past, LQCD has, therefore, been carried out at supercomputer centers. Compared to this project's computing facilities, the costs at supercomputer centers are two to five times greater to deliver the same amount of dedicated lattice QCD computing.

As a conservative estimate of the benefits resulting from the $3M investment in FY10, we use only a four-fold ROI and thus a benefit of $12M over the 4 year lifetime of the systems.  This $12M figure will be used in the BY10 Exhibit 300 submission to the OMB in the Alternatives Analysis section, with the FY10 lifecycle costs cited in the Alternatives sections above as the risk adjusted cost.